

Imperial College London

IMPERIAL COLLEGE LONDON

DEPARTMENT OF LIFE SCIENCE

Genomic Architecture of Parallel Ecological Divergence in *Littorina saxatilis*

Author:

Rui Zhang

Supervisor:

Matteo Fumagalli

m.fumagalli@imperial.ac.uk

CID:

01907894

Co-supervisor:

Francesca Raffini

f.raffini@sheffield.ac.uk

Email:

rui.zhang20@imperial.ac.uk

Roger K. Butlin

r.k.butlin@sheffield.ac.uk

Date: 8/26/2021

A THESIS SUBMITTED FOR THE PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE
DEGREE OF MASTER OF SCIENCE AT IMPERIAL COLLEGE LONDON

SUBMITTED FOR THE MSc IN COMPUTATIONAL METHODS IN ECOLOGY AND EVOLUTION

Declaration

I declare that this thesis is solely composed by my own work. The data was provided by co-supervisor Francesca Raffini and Roger K. Butlin from The University of Sheffield. The preliminary data processing was conducted by co-supervisor Francesca Raffini. I only executed data analysis and results visualization. Supervisor Matteo Fumagalli provided a workshop on how to use ANGSD v.0.934. Before August 2021, this thesis has not been published in any journal.

Signature:

A handwritten signature in brown ink that reads "Rui Zhang". The signature is written in a cursive style with a clear distinction between the two characters.

Date: 8/26/2021

Contents

Abstract	4
Keywords	4
1 Introduction	5
2 Method	8
2.1 PCA	8
2.2 Admixture analysis	8
2.3 Folded one dimensional and two dimensional site frequency spectrum	8
2.4 Population genetic differentiation and nucleotide diversity	9
2.5 Linkage disequilibrium heat map and decay analysis	9
3 Results	10
3.1 PCA	10
3.2 Admixture analysis	10
3.3 1d SFS	13
3.4 2d SFS	14
3.5 Summary statistics	15
3.6 Linkage disequilibrium	18
4 Discussion	20
4.1 PCA and admixture analysis	20
4.2 Summary statistics of differentiation	20
4.2.1 SFS	20
4.2.2 Population genetic differentiation	20
4.2.3 Limitation of FST	21
4.3 Linkage disequilibrium	22
4.4 Other possible mechanisms of parallel adaptation	23
4.5 Future direction	23
5 Conclusion	24
Code availability	24
Acknowledgement	25
Supplementary information	31
A Data collection and processing	31
A.1 Sampling	31
A.2 DNA sequencing and processing	31

B Table and plots	32
B.1 Mean summary statistics	32
B.2 PBS	33
B.3 Pi	34
B.4 LD decay	35
B.5 LD heat map	37

¹ Abstract

² Parallel adaptation of *Littorina saxatilis* in heterogeneous environments has been researched by
³ many scientists all around the world, since it is a good subject to explore sympatric speciation
⁴ and natural selection. This species has evolved two ecotypes adapted to micro-habitats: crab
⁵ ecotype and wave ecotype. The aim of this research is to find underlying mechanisms of parallel
⁶ ecological divergence with gene flow and the evolution of reproductive isolation with inter-
⁷ breeding. Low coverage whole genome sequencing and related data analysis tools were used on
⁸ 73 samples of snails and the research on hybrids is the emphasis.

⁹

¹⁰ In the Spanish system of *Littorina saxatilis*, snails are divided into two distinct genetic clusters
¹¹ instead of a cline in Sweden. The crab group and hybrid group are genetically closer and
¹² they compose the first cluster, a cline. The wave group forms the second homogeneous cluster.
¹³ Also, there is an obvious directional introgression from wave into crab at a low rate which
¹⁴ destroys high Linkage Disequilibrium (LD) in crab ecotype. Differentiation analysis indicates
¹⁵ that most loci are polymorphic and not fixed suggesting that loci are under both divergent
¹⁶ selection and balancing selection. Usually, chromosome inversion candidates show higher FST
¹⁷ and LD compared with whole chromosome, while PBS and genetic diversity (π) are more diverse.
¹⁸ The crab group has the highest LD and slowest LD decay rate which hints at stronger selection
¹⁹ and more polymorphic inversions. It is hard to find new chromosome inversions just from LD
²⁰ heat map, but some chromosomes demonstrate earlier points dispersal in LD decay suggesting
²¹ short inversions. In conclusion, this research found some new chromosome inversions candidates
²² to be verified and tried to explain the process of ecological divergence after the phase of a cline
²³ in terms of genetic structure and selection.

²⁴ Keywords

²⁵ Whole genome sequencing, parallel evolution, local adaptive, genomic architecture, gene flow,
²⁶ reproductive isolation

27 1 Introduction

28 How speciation occurs is a fundamental biological problem. There are three traditional geo-
29 graphic modes of speciation: allopatric speciation, parapatric speciation, and sympatric specia-
30 tion (Fitzpatrick et al. 2008). Sympatric speciation is fascinating due to its difficulty in theory
31 explanation and definition. Parallel ecological divergence is one of the most interesting subjects
32 of sympatric speciation. A monomorphic population could split into two coexisting phenotypic
33 clusters (ecotypes) under directional selection of ecological environments explained by the theory
34 of adaptive dynamics (Dieckmann & Doebeli 1999). It is well established that divergent natural
35 selection in heterogeneous ecological environments is likely to be an impetus for local adaptation
36 and subsequent reproductive isolation even speciation (Schluter 2009). Repeated occurrence of
37 ecological divergence in discrete populations in contrasting habitats gives an opportunity to ex-
38 plore speciation mechanisms. Moreover, how reproductive isolation evolves with interbreeding,
39 which means the presence of gene flow impeding speciation, is more informative to understand
40 initial steps of speciation than the study on already separated species (Kirkpatrick & Ravigné
41 2002).

42

43 *Littorina saxatilis* is a common rocky-shore gastropod with short life and high lifetime dispersal
44 due to ovoviparity (Reid 1996). This species has a strong adaptation ability. Even small
45 patches of local habitat could promote a distinct ecotype under biotic or abiotic pressure (Joh-
46 hannesson 2016). It has evolved two ecotypes under the selection of wave exposure and crab
47 predation in many locations such as northwest Spain, west coast of Sweden, and northeast coast
48 of England (Butlin et al. 2008). Wave ecotype is smaller and has a thinner shell with a larger
49 aperture, while crab ecotype is bigger and has a thicker shell with a smaller aperture (Johannes-
50 son et al. 2010). The morphology and behavior differences between two ecotypes are inheritable
51 although with a minor phenotypic plasticity (Galindo et al. 2009). *Littorina saxatilis* has devel-
52 oped partially reproductive isolation in the face of hybridization (Grahame et al. 2006). This
53 proposition is supported by that gene flow in contact zones was 10-30 % of gene flow within
54 ecotypes and that the genetic relationship between crab ecotype and wave ecotype in the same
55 location is closer than the genetic relationship between the same ecotype in different locations
56 (Panova et al. 2006). Also, Beaumont (Beaumont 2010) got strong support for parallel and local
57 divergence in this species instead of old allopatric divergence of ecotypes followed by secondary
58 overlap and gene flow using approximate Bayesian computation (ABC) approach. Therefore,
59 it is a preferable model system to study the underlying mechanisms, particularly genomic ar-
60 chitecture variation, of parallel adaptive divergence with the existence of gene flow, which is
61 sympatric or parapatric speciation instead of allopatric speciation.

62

63 There are many studies focusing on the history of *Littorina saxatilis* dispersal, colonization,
64 formation of ecotypes, and evolution of reproductive barriers between two ecotypes under gene
65 flow. This species is inferred to survive from the last glacial period in a northern latitudes
66 refugia based on phylogeographic data (Panova et al. 2011). Later colonization seemed to be
67 achieved by rafting of single females whose brood pouch carried hundreds of embryos sired by

even over 20 males. This feature allows a single female to be a founder group with diversified genetic variation and a new population is established rapidly once releasing embryos (Rafajlovic et al. 2013). Crab is regarded as the driving force of crab ecotype formation. After snails' colonization, the predator of snails arrived at the same location later due to their higher minimum temperature requirement compared with snails. In the ABC model established by Beaumont (Beaumont 2010), ecotypes separation is a relatively recent incident. The ecotype separation time is estimated at only around 10 % of local population age. Furthermore, the formation of ecotypes occurs instantaneously (<1000 generations) rather than gradually. Although ecotypes have formed, individuals in two ecotypes are still able to interbreed with each other even from different locations (Hollander et al. 2005). The survival rate of hybrids in the contact zone is higher than parental ecotypes or on the same level which means hybrid superiority exists (Rolan-Alvarez et al. 1997). It promotes the first step of speciation, ecotype formation because hybrids are easy to back-cross with parental species. However, it impedes the completion of speciation due to continuous gene flow.

82

Gene flow between crab and wave ecotypes never disappear although in the contact zone gene flow is impeded and less than gene flow within ecotypes (Panova et al. 2006). How *Littorina saxatilis* maintain divergence in the face of gene flow which is regarded as genetic recombination counteracting divergence? Strong enough selection pressure could overcome this kind of homogenizing effects and genomic architectures might resist gene flow by impeding gene recombination (Smadja & Butlin 2011). Hybrid zone analysis has inferred patterns of selection in space of *Littorina saxatilis* that crabs represent a strong selection pressure during ecotypes separation (Westram et al. 2018). The genetic basis of parallel adaptation is supposed to depend on underlying genomic architecture (Yeaman 2013) such as the number, size and additivity of loci and nonrandom arrangements order in the genome. Chromosome inversion, a kind of chromosome rearrangement, has been researched widely. It is considered to be able to suppress recombination and serve as reservoirs for adaptive standing variation which might be the reason for fast parallel adaptation (Morales et al. 2019). In addition, assortative mating of ecotypes is obvious in both field and laboratory experiments. Individuals are more likely to mate with the same ecotype individuals by means of following trails and longer mating time (Hollander et al. 2005). In brief, the partial speciation of this species is considered to be evolved by divergent selection of heterogeneous environments, habitat choice, assortative mating, and genomic background (chromosome rearrangement).

101

In our research site, Spanish, ecotypes are distributed over vertical shore gradients. Hybrids are a minority compared with parental ecotypes in the intermediate habitat (Johannesson et al. 1993) with isolation indexes of 0.5–0.9 in Spain tested from mate choice experiments (Johannesson et al. 1995). The Spanish system of *Littorina saxatilis* ecological divergence is found much older than Swedish and British systems based on mitochondrial DNA lineages (Panova et al. 2011). Furthermore, Morales (Morales et al. 2019) found that the comparisons between crab ecotype and wave ecotype showed more significant divergence in Spain than in Sweden. However, the

109 comparisons only focus on crab ecotype and wave ecotype snails themselves and ignore hybrids.
110 Hybrid of two ecotypes could offer more information than distinct population comparisons. For
111 instance, cline analysis and hybrid zone analysis in Sweden (Westram et al. 2018) were used
112 to detect non-neutral SNP (single nucleotide polymorphism) and find a genotype-phenotype-
113 environment association. Therefore, this project will take hybrid whole genome sequencing into
114 account.

115

116 Then what kind of sequencing method and subsequent data analysis should be used in this
117 research? Low coverage whole genome sequencing (lc WGS) is regarded as a cost-efficient ap-
118 proach to catch low frequency variation in many individual samples of population (Gilly et al.
119 2018). The development of calling algorithms makes it possible to call SNP and estimate allele
120 frequency accurately from low depth sequencing data. Analysis of next generation sequencing
121 data (ANGSD) is a software designed for next generation sequencing data with data analysis
122 methods of taking genotype uncertainty into account instead of calling genotypes directly (Ko-
123 rneliussen et al. 2014). It is especially suitable and useful for low and medium depth data.
124 Therefore, it is a good sequencing data analysis tool for lc WGS data of population genomics.
125 ANGSD has been used to explore underlying genetic mechanisms of parallel phenotypic evolu-
126 tion, local adaptation, history of evolution, and so on (Wilder et al. 2020) (Therkildsen et al.
127 2019) based on population genetic sequencing data. Furthermore, linkage disequilibrium is a
128 good indicator to explore population genetic history including selection, domestication, chromo-
129 some rearrangement, and so on. Also, LD decay analysis could reveal population recombination
130 and selection history (Zhang et al. 2018). All of these measures help us to explore population
131 structure, natural selection signature, and history of genetic changes.

132

133 In this research, I expect to find a genetic continuous cline like the Sweden system from principal
134 components analysis (PCA), special population genetic differentiation, and nucleotide density
135 pattern in chromosomes, detect new inversions candidates and infer parallel adaptive evolution
136 history of *Littorina saxatilis* in Spain. Combining all above, I try to explain the underlying
137 mechanism of parallel ecological divergence and the process of speciation facing gene flow.

¹³⁸ **2 Method**

¹³⁹ Samples collection, DNA sequencing, and processing were conducted by staff at the University
¹⁴⁰ of Sheffield. Also, details of these processes were provided by co-supervisor Francesca and Roger
¹⁴¹ and described in the supplementary information.

¹⁴² **2.1 PCA**

¹⁴³ After getting bam files, I used software ANGSD v.0.934 (Korneliussen et al. 2014) (<http://www.popgen.dk/angsd/index.php/ANGSD>) to estimate imputed genotype probabilities from mapped
¹⁴⁴ reads. ANGSD is a software designed for analyzing next generation sequencing data, especially
¹⁴⁵ low and medium depth data due to its feature of taking genotype uncertainty into account. In
¹⁴⁶ order to explore the population structure of Littorina saxatilis in Spain, I performed principal
¹⁴⁷ components analysis (PCA) and admixture analysis. At first, I used ANGSD to perform SNP
¹⁴⁸ calling, estimate genotype likelihood in beagle format and get a list of SNP positions as input
¹⁴⁹ files to estimate covariance matrix. Both ANGSD and PCAngsd could infer covariance matrix
¹⁵⁰ of Littorina saxatilis 73 samples. Here I chose PCAngsd v.1.02 (Meisner & Albrechtsen 2018)
¹⁵¹ (<http://www.popgen.dk/software/index.php/PCAngsd>) and then used R v.4.0.3 to perform
¹⁵² PCA using eigen function and extract eigenvectors to plot PC1 vs PC2. I also used individuals'
¹⁵³ distance along crab-hybrid-wave axis data as color information to show population structure
¹⁵⁴ along crab-hybrid-wave axis.
¹⁵⁵

¹⁵⁶ **2.2 Admixture analysis**

¹⁵⁷ Another method to interpret population structure is admixture analysis which showed genome-
¹⁵⁸ wide admixture proportions for every individual. ngsAdmix (<http://www.popgen.dk/software/index.php/NgsAdmix>) (Skotte et al. 2013) used beagle format genotype likelihood file to infer ad-
¹⁵⁹ mixture proportions of ancestry clusters for every individual. Then I used R to draw admixture
¹⁶⁰ proportions bar plot along crab-hybrid-wave axis. After getting the admixture analysis result,
¹⁶¹ I could classify 73 individuals into 3 groups: crab, hybrid and wave according to individual
¹⁶² ancestry clusters contribution. Given ecotype information, I could draw a new PCA plot in 3
¹⁶³ colors showing population structures of every ecotype population.
¹⁶⁴

¹⁶⁵ **2.3 Folded one dimensional and two dimensional site frequency spectrum**

¹⁶⁶ In order to detect selection signatures on genome and understand patterns of divergence and
¹⁶⁷ differentiation across genome, I calculated several summary statistics from low-depth NGS data.
¹⁶⁸ The first step was estimating sample allele frequencies (SAF) posterior probabilities by ANGSD
¹⁶⁹ for each population separately. Then I used the program realSFS of ANGSD to calculate Site
¹⁷⁰ Frequency Spectrum (SFS) which records the proportions of sites at different allele frequencies.
¹⁷¹ Here I calculated folded SFS without outgroup species defining ancestral state. After getting
¹⁷² one dimensional folded SFS, I drew a bar plot for every population removing the first value of 1D
¹⁷³ folded SFS which represented the expected number of sites with derived allele frequency equal
¹⁷⁴ to 0 due to its relatively too big value to guarantee other values clear interpretation. The next

175 step was estimating joint SFS between 2 populations (2D folded SFS) which is useful to infer
176 the divergence process of populations and estimate population genetic differentiation as prior
177 information. As before calculating 1D folded SFS, I used program realSFS to infer 2D folded
178 SFS between every pair of 3 ecotype populations. Similarly, I could use R to illustrate these
179 flatten matrixes in heat map format.

180 2.4 Population genetic differentiation and nucleotide diversity

181 Here I chose fixation index (FST) (Hudson et al. 1992) and population branch statistic (PBS)
182 as indicators of allele frequency differentiation. ANGSD helps us calculate FST and PBS from
183 sample allele frequencies likelihoods (.saf files), avoiding genotype calling. The first step was
184 computing per-site FST indexes by realSFS using SAF files and 2D folded SFS as input. The
185 output file of this stage is (a) and (a+b) values calculated in the way of Reynolds(Reynolds
186 et al. 1983) for three FST comparisons for every SNP site. I could use R to calculate FST value
187 for every SNP and every pairwise of 3 populations and then drew FST in every chromosome
188 and every pairwise populations in the order of position. The next step was performing a sliding-
189 window analysis with setting window size as 10kb and step size as 1kb by realSFS. Similarly I
190 drew window PBS values in every chromosome and population.

191 Furthermore, in order to know whether allele frequency differentiation increase is related to the
192 change of nucleotide diversity (Korneliussen et al. 2013), I used realSFS saf2theta function and
193 program thetaStat of ANGSD to calculate thetas for each site and in windows taking SAF and
194 SFS as input files and prior information. Pairwise theta estimation (t_P) divided by numbers of
195 sites in the window (n_{Sites}) can be used to estimate window-based pairwise nucleotide diversity
196 (π).

198 2.5 Linkage disequilibrium heat map and decay analysis

199 Linkage disequilibrium is a useful indicator to illustrate the correlation of pairwise loci and
200 identify inversions. Here are various different measures of LD such as pairwise r^2 , Pearson
201 estimation of r^2 , D from EM algorithm and so on. Here I choose pairwise r^2 (coefficient of
202 correlation) because it is very useful when analyzing biallelic markers such as SNPs and inde-
203 pendent of sample size (Devlin and Risch 1995). I subset beagle files and maf files of every
204 chromosome at first, and then sample these SNP sites to guarantee the total SNP number of
205 every chromosome is around 10000 as instruction of ngsLD software. ngsLD v.1.1.1 (Fox et al.
206 2019)(<https://github.com/fgvieira/ngsLD>) is a software designed for NGS data taking the un-
207 certainty of genotype's assignation into account. After using this software calculating LD among
208 SNPs in every chromosome, I used R to draw LD heat map in whole chromosome and partial
209 chromosome inversion candidates. Then I used a modified R script offered by ngsLD to analyze
210 LD decay and draw LD decay line for every chromosome and population.

211 **3 Results**

212 **3.1 PCA**

213 The covariance matrix of principal components analysis (PCA) could be estimated by single-
214 read sampling in ANGSD. With proper parameters settings, there are 12,817,872 variant sites
215 (SNP) called from 73 individuals' bam files. Here I used PCAngsd to take genotype likelihoods
216 in beagle format as input and then infer covariance matrix which is used in Figure 1.

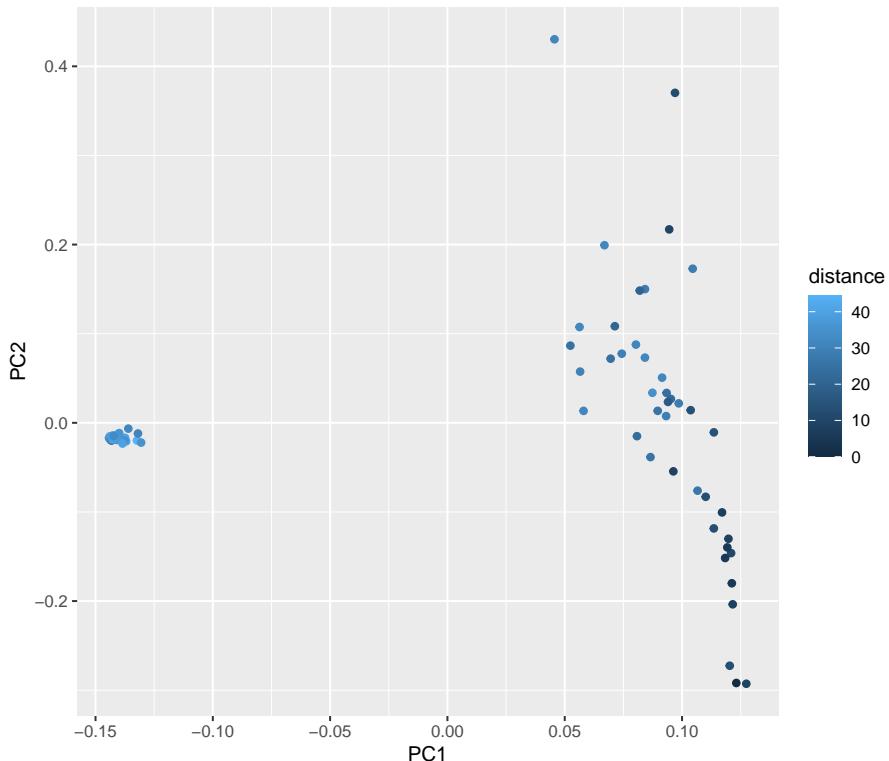


Figure 1: **PCA of 73 Littorina saxatilis individuals**

Dark blue presents crab ecotype and light blue presents wave ecotype. The more dark means the sampling individual is closer to crab ecotype along the crab-hybrid-wave axis. V2 is the first principal component and V3 is the second principal component.

217 The PCA result was quite different from the Sweden system. There are many individuals points
218 squeezing together in a little circle region and the color of these individuals is light blue which
219 means these individuals are likely to be wave ecotype snails. On the other side of Figure 1,
220 points are darker and comparatively scattering. Furthermore, deep dark blue points are usually
221 lower than medium blue points. So far, I could only know that there are 2 genetic distinct
222 clusters along crab-hybrid-wave axis.

223 **3.2 Admixture analysis**

224 In order to infer every individual's genome-wide admixture proportions, I used ngsAdmix to
225 estimate admixture proportions from genotype likelihoods and ancestry clusters number setting.

226

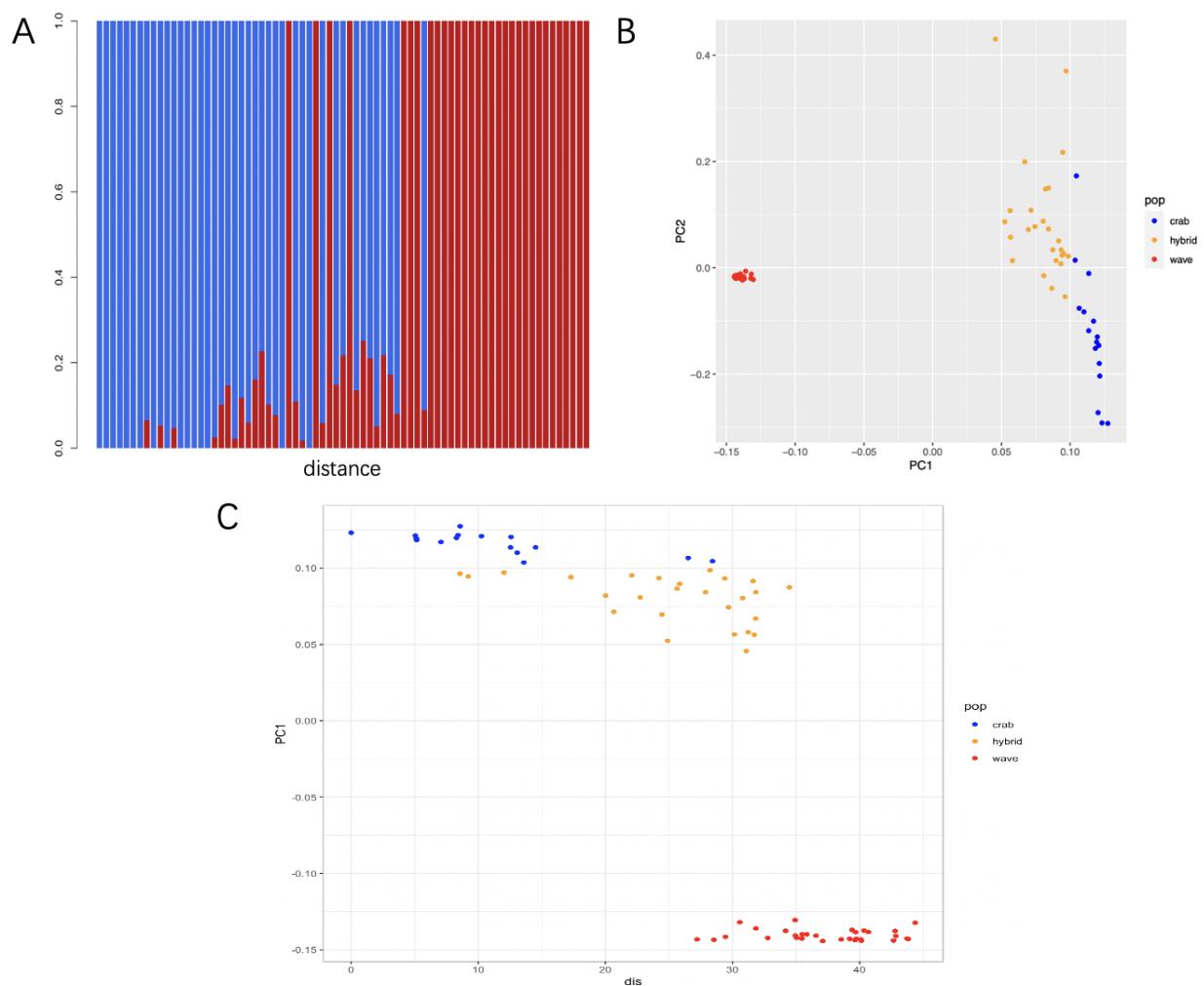


Figure 2: Admixture analysis of 2 ancestry clusters and new PCA based on admixture analysis

- A. The figure showed individuals admixture analysis results along crab-hybrid-wave axis from crab to wave ecotype. Blue represents crab ecotype and red represents wave ecotype.
- B. Based on admixture analysis, individuals could be divided into 3 groups: wave, hybrid and crab. Given the group information, PC1 vs PC2 of principal components analysis is shown in 3 colors.
- C. Based on admixture analysis, 3 groups named wave, hybrid and crab could be plotted in a PCA between distance of crab-hybrid-wave axis and PC1.

227 As Figure 2A showed, right individuals are mostly wave ecotypes while left individuals are mostly
228 crab ecotypes which corresponds with distance data along crab-hybrid-wave axis. However, in
229 the middle part of Figure 2A, there are crab ecotype, wave ecotype and hybrid individuals,
230 which means in the contact zone, 2 ecotypes and their hybrid offspring are mixed. Furthermore,
231 hybrid individuals can also enter pure wave and pure crab ecotype zone. From the result of
232 admixture analysis, I could divide individuals into 3 groups: crab ecotype, wave ecotype and
233 hybrid with 16, 31 and 26 individuals separately. Hybrids include all individuals with admixture
234 proportions that are not 0 or 100%. When group information is considered, PCA plot could be
235 redrawn as Figure 2B. All of wave ecotype individual points squeezed together just as Figure 1.
236 Crab ecotype individuals and hybrid individuals are all on the other side, far away from wave
237 ecotype, which means crab group and hybrid group are much more close compared with wave
238 group. Meanwhile, the hybrid group is higher in the second principal component than crab
239 group, although it is difficult to distinguish hybrid group and crab group only using the first
240 two principal components.

241

242 In Figure 2C, all wave individuals are around -0.14 in PC1, all crab individuals are higher
243 than 0.10 while hybrid individuals are lower than 0.10 in PC1. However, only using the line of
244 $PC1=0.10$ is not sufficient to split the higher cluster into crab and hybrid two isolated popula-
245 tions. Therefore, I would like to regard that there are 2 distinct groups, one is a cline within
246 crab and hybrid group, the other is wave group.

²⁴⁷ 3.3 1d SFS

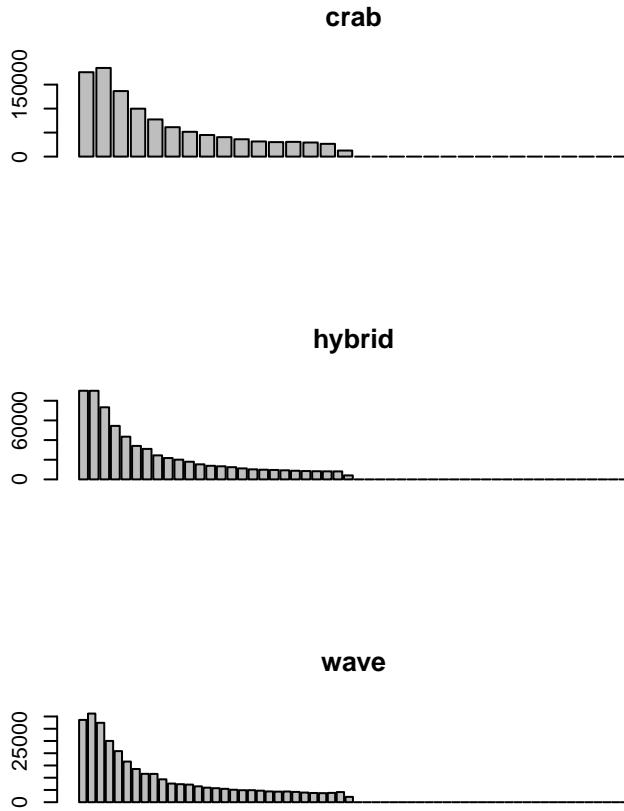


Figure 3: **One-dimensional folded Site Frequency Spectrum (SFS)**

One-dimensional folded SFS of each group is shown in 3 bar plots from allele frequency 1 to the number of individuals in the group

²⁴⁸ Site frequency spectrum records the number of sites at different allele frequencies. One-dimensional
²⁴⁹ folded SFS only contains allele frequencies in one group. In Figure 3, sites will be recorded in
²⁵⁰ SFS only if this site is caught in every individual of the group, also I have removed all of the sites
²⁵¹ without alleles due to its comparatively huge value. As the bar plots showed, there are fewer
²⁵² sites of allele frequency 1 than allele frequency 2 in crab and wave group. It is unexpected but it
²⁵³ might be caused by distortion effect on SFS of polymorphic inversions and SNP filter conditions.
²⁵⁴ The sites number of crab group is bigger than hybrid and wave group. It could be explained
²⁵⁵ by smaller population size of crab group which means it is easier for crab group to detect the
²⁵⁶ same site in every individual. Furthermore, after calculating the proportion of SNP sites in
²⁵⁷ every group, crab group SNP proportion is 0.0193, wave group SNP proportion is 0.0185, hybrid
²⁵⁸ group SNP proportion is 0.02266. The hybrid group has the highest SNP proportion maybe
²⁵⁹ because the hybrid group contains SNP both in crab group and wave group.

260 3.4 2d SFS

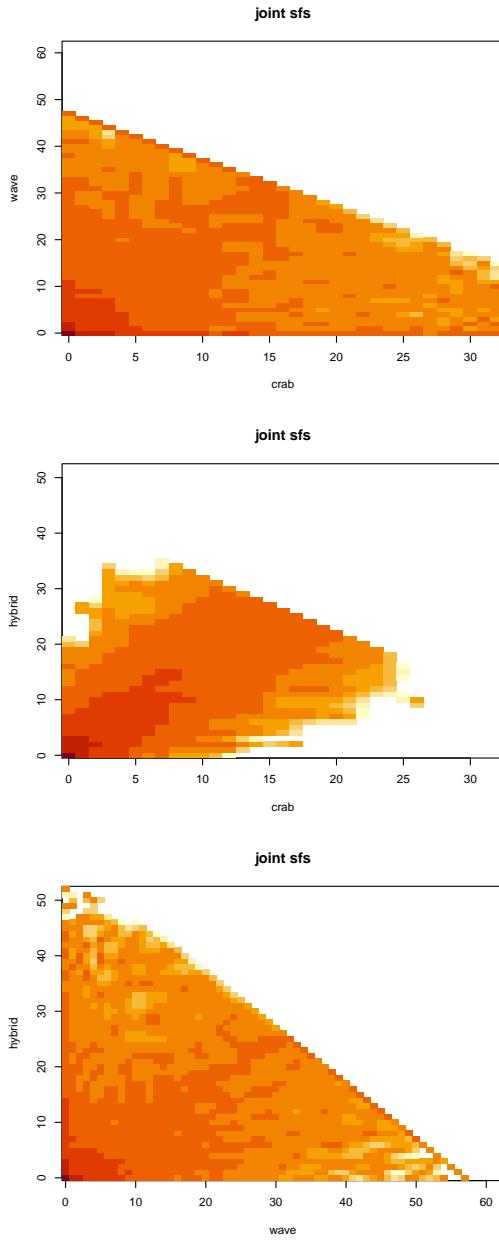


Figure 4: **Two-dimensional folded Site Frequency Spectrum (SFS)**

Two-dimensional folded SFS of each pair of groups are shown in 3 heat plots from allele frequency 0. Darker color means more sites and lighter color means less sites.

261 Joint SFS between 2 populations could be used to infer the divergence process. From Figure 4,
 262 it is clear that crab and hybrid groups are closer because there are darker color squares on the
 263 diagonal compared with the other 2 heat plots. Usually colored squares of folded 2D-SFS are all
 264 under one of the diagonals. However in Figure 4, 2D folded SFS didn't meet the expectation.
 265 This situation might be resulted from different population sizes. 2D-SFS could also be used
 266 as prior information for estimating FST and PBS (population branch statistic) which will be
 267 discussed later.

268 3.5 Summary statistics

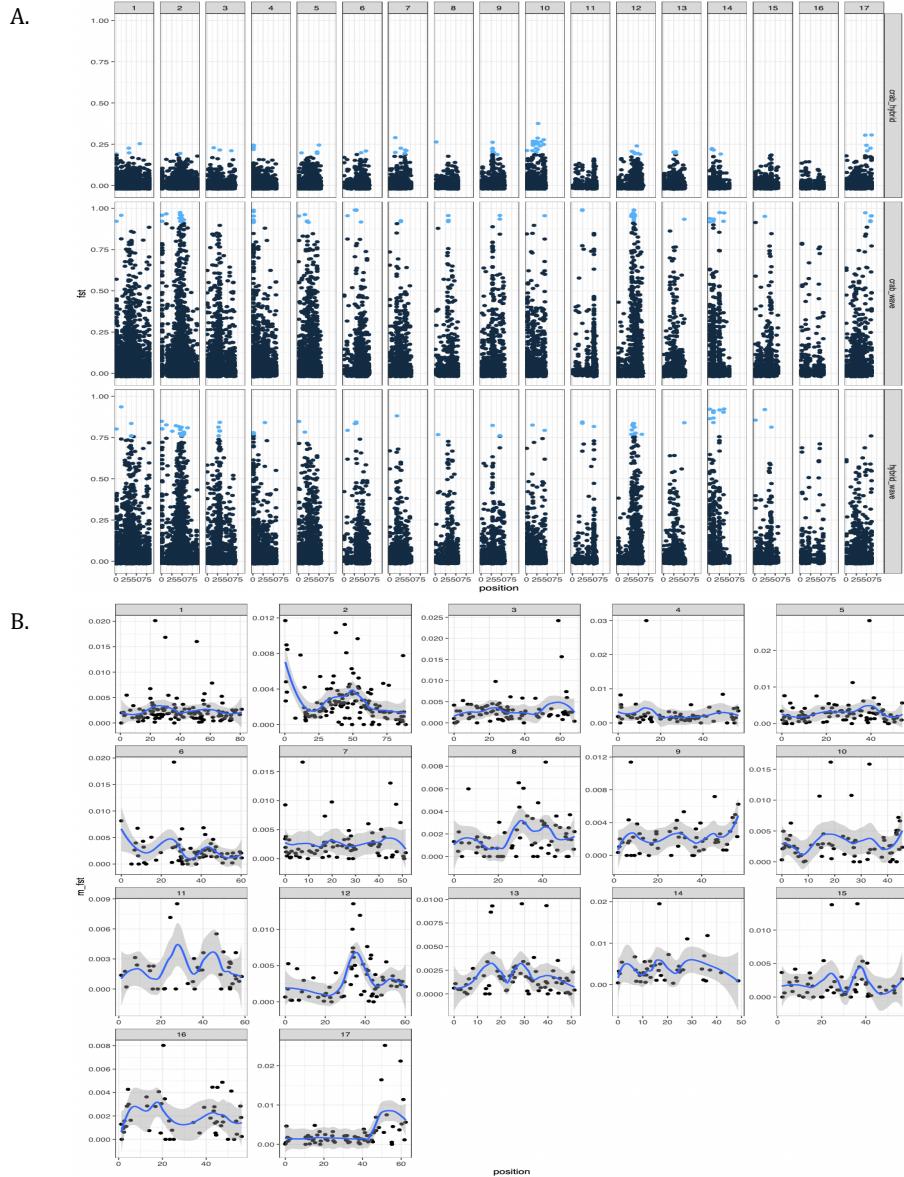


Figure 5: FST value among 17 chromosome and three pairwise populations

A. FST value scatter plot of every chromosome and pairwise populations with blue points representing outliers.

B. Mean FST value between crab and wave ecotype in every position on chromosome along with smooth lines showing trend.

269 There are some big blocks showing high FST value, high PBS value and low pi value, which
 270 are all signatures of positive selection, especially in LG 1, 2, 4, 6, 9, 12, 14, 17. From Figure
 271 5A, I can find that population genetic differentiation of crab and wave is the largest. Also, it is
 272 clear that hybrid and crab are closer due to their low FST value in all chromosomes. Figure 5B
 273 exhibits mean FST value of crab and wave population showing general high FST value in the
 274 specific position such as crest part of LG2, 6, 12, 17. Population branch statistics of hybrid is
 275 the lowest which is not difficult to understand because hybrids contain genetic components from

276 crab and wave ecotype. PBS of wave is the highest which is another proof of closer relationship
277 between crab and hybrid. Pairwise nucleotide in Figure 2 in supplementary information can not
278 offer too much effective information. The pi value looks nearly all the same in every chromosome
279 and every ecotype.

280

281 Faria et al. (Faria et al. 2019) found some chromosome inversion candidates as Figure 6A
282 shown. In these regions, FST and PBS are usually relatively higher than other regions, while pi
283 are usually relatively lower than other regions. In order to interpret differences between these
284 chromosomal rearrangement candidates and other regions, I can also calculate mean value of
285 FST, PBS, pi and linkage disequilibrium (LD, which will be discussed later) in these regions
286 and whole chromosome. Table S1 in the supplementary information gives the exact number of
287 these statistics values which is more helpful to some extent. Also, a scatter plot Figure 6B has
288 been drawn to interpret differences between candidates and whole chromosome more intuitively.
289 From Figure 6B, generally FST and LD values of chromosome inversion candidate regions are
290 higher than whole chromosome. However, PBS and pi value of chromosome inversion candidates
291 region is sometimes lower and sometimes higher than whole chromosome. Statistics of LG1.2,
292 2.1, 4.1, 6.1, 6.2, 9.1, 14.3 meet the expected result of higher PBS and lower pi in inversion
293 candidates regions. LG1.1, 7.1, 7.2, 17.1 show higher PBS and higher pi while LG 10.1, 11.1,
294 12.1, 12.2, 14.1, 14.2 show lower PBS and higher pi.

295

296 However, there are also some chromosome regions showing high FST and PBS and low pi, but
297 are not chromosomal rearrangement candidates, such as the middle part of LG 2, 8, 11 and 15.
298 They will be discussed later.

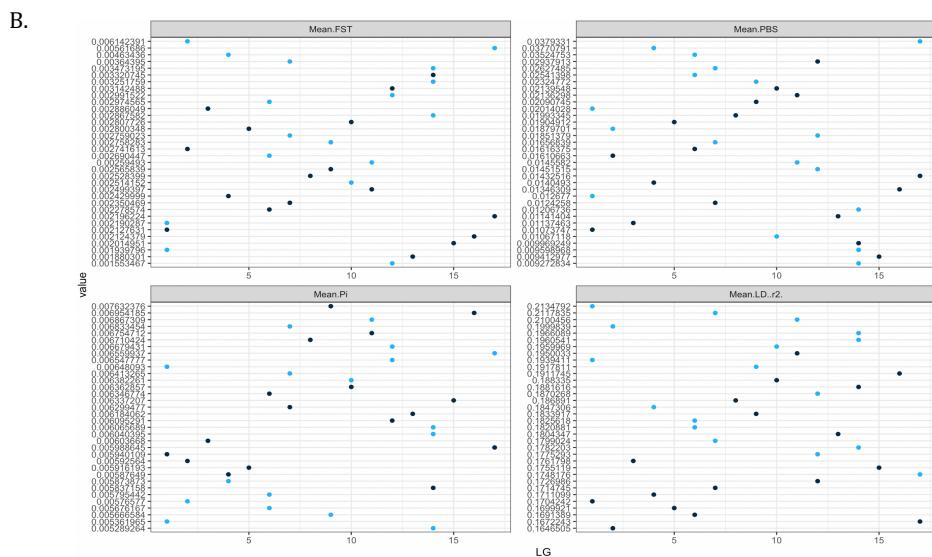
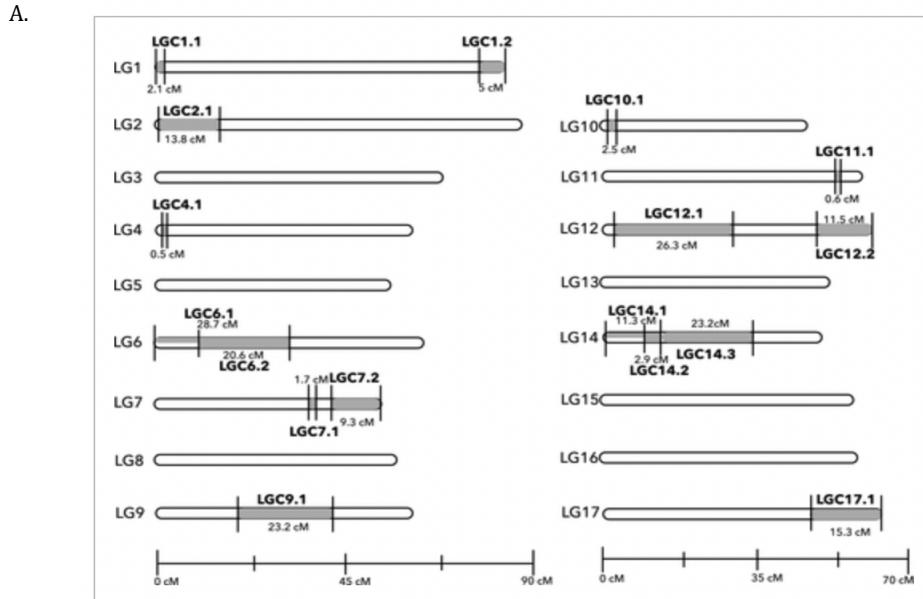


Figure 6: Mean statistics of inversion candidates and whole chromosome

A. Chromosomal rearrangement candidates inferred by Faria et al.(Faria et al. 2019)

B. Blue represents inversion candidates mean statistics while black represents whole chromosome mean statistics. The statistics here include mean FST of crab and wave ecotype, mean PBS of crab ecotype, mean pi of crab ecotype and mean LD of crab ecotype.

299 **3.6 Linkage disequilibrium**

300 Linkage disequilibrium estimates the nonrandom association of alleles at different loci. Natural
 301 selection, mutation, non-random mating, gene flow and population size can all affect LD. Here
 302 I choose pairwise r² (coefficient of correlation) because it is very useful when analyzing biallelic
 303 markers such as SNPs and independent of sample size (Devlin and Risch 1995). All of the LD
 304 heat map and LD decay analysis regression plots are included in Figure 3 and 4 of supplementary
 305 information. Here Figure 7 only gives an example on chromosome 14. In 17 chromosomes, the
 306 crab population always has the highest LD value. Hybrid usually has medium LD value and
 307 wave has the lowest LD value.

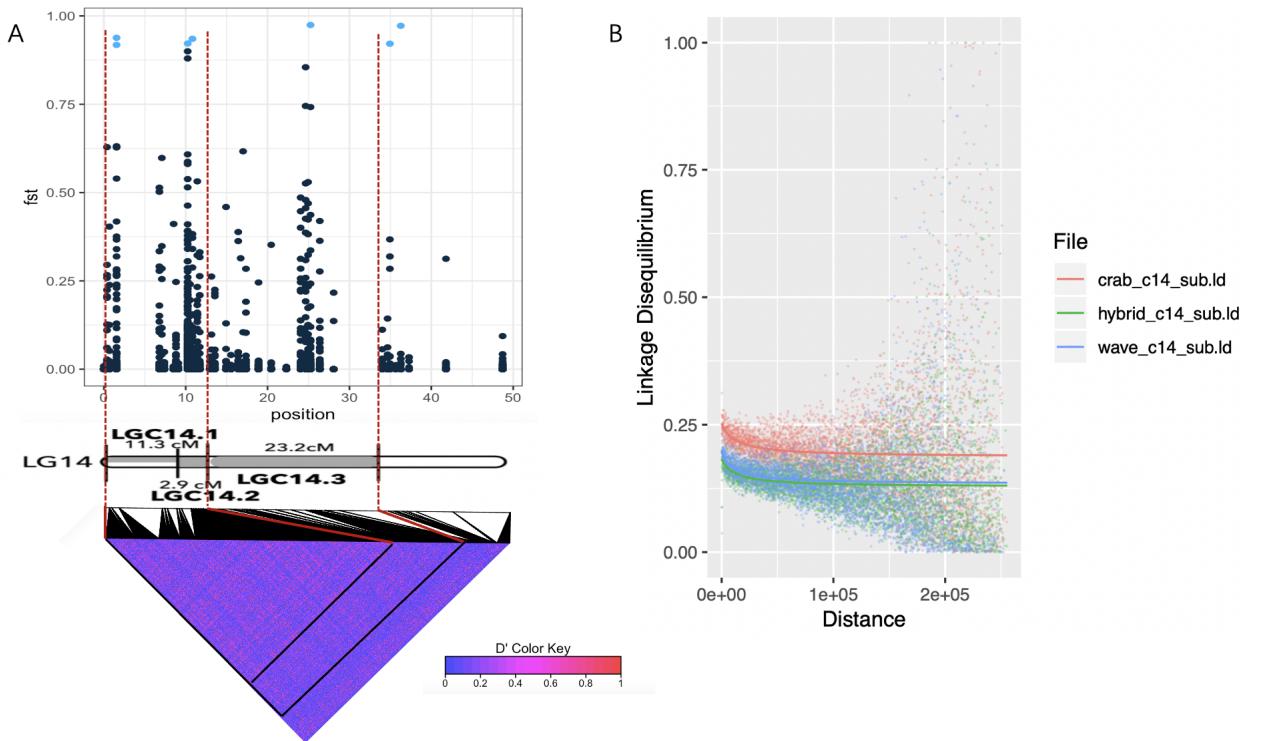


Figure 7: **LD heat map and LD decay of chromosome 14**

A. FST value with outliers mapping to chromosome inversion candidates and LD heatmap of crab population on chromosome 14.

B. LD decay analysis of chromosome 14

308 LD heat map among the whole chromosome is hard to get useful information because usually
 309 all regions look the same. However in LG 14, the heat map between the first two red dotted
 310 lines has a clear higher LD value block which maps to LG14.1 and LG14.2. The region between
 311 the second and the third red dotted lines is also more pink than white region on chromosome
 312 14. Partial LD heat map in chromosome inversion candidates is more clearly showing higher
 313 overall LD value compared with other regions. Exact mean LD values are shown in Table S1
 314 that mean LD values of chromosomal rearrangement candidates are higher than other regions.

315

316 For LD decay analysis, every single point is mean LD value of a 50bp size window. The fitting
 317 line decreases quickly approaching a horizontal line. As the distance increases, points become

³¹⁸ more dispersed and some chromosomes disperse more early such as LG1 and LG4. Crab usually
³¹⁹ has the slowest LD decay rate in 17 chromosomes while hybrid ecotype population LD decay
³²⁰ rate is usually the fastest, which might be resulted from higher population genetic diversity
³²¹ gained from crab and wave ecotype population.

322 **4 Discussion**

323 **4.1 PCA and admixture analysis**

324 There are 2 distinct clusters in PCA in the Spain system instead of a cline which might be
325 resulted from more time to accumulate divergence. The first group is more homogeneous and
326 regarded as wave ecotype. The second group is more heterogeneous and regarded as a cline
327 within crab and hybrid group. Furthermore, from the admixture analysis, I could find obvious
328 directional introgression from wave into crab at a low rate.

329

330 In Sweden admixture analysis conducted by Sheffield University staff, there were some indi-
331 viduals showing higher than 50 % of wave contribution, but in the older system of Spain, the
332 situation is different. A possible speculation is that, in the process of speciation, hybrids are
333 more likely to mate with crab (or to say bigger individuals), then hybrid population in PCA
334 become closer and closer to crab population, ancestry proportion of wave decrease and propor-
335 tion of crab increase. In addition to assortative mating, the possible reason for one direction of
336 introgression could also be asymmetrical selection, or confusion of adaptation to the lower shore
337 within crab group or a mix of these.

338 **4.2 Summary statistics of differentiation**

339 Here I expect to find special population genetic differentiation and nucleotide density pattern
340 in chromosomes, detect new inversions candidates and infer parallel adaptive evolution history.

341 **4.2.1 SFS**

342 The 1d SFS seems to be far from the neutral expectation. It could not be caused by introgression
343 since it is similar in all groups. The SFS might be distorted by SNPs in polymorphic inversions.
344 I only used sites detected in all individuals to guarantee the accuracy of gene frequency, but
345 this might also distort the SFS. For instance, if repeat regions are more likely to fall into this
346 class, the SFS will be inaccurate. From 2d SFS, the only significant information is that crab
347 and hybrid populations are closer which is consistent with PCA result.

348 **4.2.2 Population genetic differentiation**

349 Most of loci FST value doesn't reach 1 which means these loci are not fixed and keep poly-
350 morphic. Rare fixed loci suggest balancing selection which tends to prevent fixation. The same
351 pattern for SNPs was found by Westram (Westram et al. 2018) and some possible explana-
352 tions were given: indirect divergent selection on SNPs linked to selected variants, selection on
353 polygenic traits or a combination of divergent selection and balancing selection that make loci
354 maintain polymorphism in one or two habitat ends.

355

356 High FST means high divergence. Regions with lots of high FST values may well contain genes
357 under divergent selection between crab and wave ecotypes although along with some noises. It is

358 hard to find these regions just from FST scatter plot due to some regions' more SNPs. Usually,
359 regions with many SNPs tend to own high FST SNPs. Therefore, I colored the top 0.01% of
360 SNPs as Figure5A and calculated mean FST per map position showed as Figure5B. There are
361 some of the same regions detected before by Morales (Morales et al. 2019) including inversion
362 candidates. There are also some regions on LG2, 8, 11 demonstrating high population genetic
363 differentiation but not regarded as inversion candidates. However, further exploration should be
364 executed to verify whether they are new inversions or just genetic speciation islands that harbor
365 loci underlying reproductive isolation, or regions of low recombination.

366

367 Generally, I expect to find high FST, high PBS, low pi and high LD regions as candidate blocks
368 for reproductive isolation. Chromosome inversion candidates are supposed to carry many loci
369 under divergent selection, so I expect to detect the same summary statistics pattern in these
370 candidates. However, I found that many candidates show higher pi, even lower FST and PBS.
371 Only LG1.2, 2.1, 4.1, 6.1, 6.2, 9.1 meet all conditions. Higher pi is possible when inversions are
372 polymorphic. The elevated diversity could be between arrangements or within an arrangement
373 at the nucleotide level due to balancing selection. For example, LG10.1, 12.1, 12.2, 14.1, 14.2
374 show lower mean FST, lower mean PBS but higher mean pi and higher mean LD. Balancing
375 selection may also result in lower FST and lower PBS, but these regions still maintain high
376 LD. Some inversion regions have similar levels of diversity to whole chromosome suggesting
377 non-polymorphic and not accumulated differentiation. Inversion regions have less diversity sug-
378 gesting these could have swept to high frequency recently.

379

380 It is worth noting that chromosome inversions can't fully explain parallelism. Morales also found
381 shared outlier loci distributing across the genome (Morales et al. 2019). This could be caused
382 by polygenic selection of multiple loci with small effect underlying parallel adaptive divergence.

383 4.2.3 Limitation of FST

384 In genomic regions surrounding barrier loci, the barrier effect initially allows a build-up of ge-
385 netic differentiation in the form of allele frequency variation between populations, typically using
386 a relative measure such as FST (Ravinet et al. 2017). However, Using FST as an indicator of
387 divergent selection has some problems.

388

389 The pattern of FST is affected by multiple factors that change throughout the genome, including
390 mutation, genetic drift, selection, demographic history, recombination, gene flow, gene density,
391 and genome structure, some of which are expected to change at different stages of speciation.
392 Polygenic makes it possible that FST of many underlying loci remain low when snails achieving
393 trait divergence(Westram et al. 2014). Background and positive selection may produce similar
394 patterns in genome scan using FST, reducing intraspecific diversity and increasing interspecific
395 gene composition(Cruickshank & Hahn 2014). Furthermore, some non-inversion candidates
396 regions with high FST may not be other currently undetected chromosomal inversions under
397 positive divergence selection, but related to other mechanisms suppressing recombination such

398 as centromere region combined with background selection or divergence hitchhiking (Ravinet
399 et al. 2017).

400 **4.3 Linkage disequilibrium**

401 From LD analysis, I could find that in 17 chromosomes, the crab population always has the
402 highest LD value. LD is likely to be high in crab for many reasons. It might be because crab
403 has a small effective population size, or crab suffered bottleneck accident in history, or crab
404 is under a relatively strong selection, or crab individuals have a stronger tend of assortative
405 mating. Genetic structure may also increase LD. Polymorphic inversions are able to suppress
406 recombination. The clines generate LD whether they are due to introgression or due to a se-
407 lective gradient. Generally, hybrid has medium LD value and wave has the lowest LD value,
408 but the difference is little. The selection of wave environment is not as strong as crab, so it is
409 reasonable that LD value of wave ecotype population is the lowest, and hybrid is medium. The
410 little LD difference between hybrid and wave groups indicates that the directional introgression
411 destroys high LD genetic structure in crab group.

412

413 LD heat map should reveal inversions and a general pattern of higher LD among neighboring
414 loci, especially among SNPs in the same contig. However, in LD heat map of 17 chromosomes in
415 *Littorina saxatilis*, there is usually no significant higher LD value block. Chromosome inversion
416 candidates LG4.1, 7.2, 9.1, 11.1, 14.1, 14.2 are relatively clear to find in whole chromosome LD
417 heat map compared with other inversion candidates. Therefore, it is more difficult to infer new
418 inversion candidates just from LD heatmap.

419

420 For LD decay analysis, crab usually has the slowest LD decay rate in 17 chromosomes which
421 means the selection effect on crab ecotype population is stronger than other populations. Se-
422 lection will decrease population genetic diversity and strengthen association between pairwise
423 loci (LD). Therefore, usually population under stronger selection has a slower LD decay rate.
424 Littorina is a high fertility and short generation interval species, therefore LD decay rate in this
425 species is generally fast. Hybrid ecotype population LD decay rate is usually the fastest which
426 might be resulted from higher population genetic diversity gained from crab and wave ecotype
427 populations.

428

429 Points disperse along the distance may be resulted from polymorphic inversions on chromo-
430 somes which makes high LD between SNPs far apart. In some chromosomes, points disperse
431 more early than others which might be because short inversions on chromosome which result in
432 LD between loci far away. The most obvious effect of inversions on LD decay is typically going
433 to be on larger scales especially between breakpoints far apart (such as 10Mb), but due to the
434 limitation of contig size, I could not find this situation.

435

436 Ravinet (Ravinet et al. 2016) found that LD clusters showed strong signals of ecotype specificity,
437 loci were largely homozygous in crab while heterozygous in wave, suggesting these loci have

438 undergone selective sweeps (the reduction or elimination of variation at sites that are physically
439 linked to a site under directional selection) in crab population. This is another proof that crab
440 group suffered stronger selection fueling divergence.

441 4.4 Other possible mechanisms of parallel adaptation

442 Morales (Morales et al. 2019) supposed that standing variation within chromosomal inversions
443 can be maintained as balanced polymorphism and fuel rapid parallel phenotypic divergence to
444 heterogeneous environments through gene flow without high sharing of genomic outliers. Bal-
445 ancing and divergent selection between habitats could maintain inversions for long periods of
446 time, resulting in the high diversity and divergence for some inversions noted above (Faria
447 et al. 2019). Furthermore, balancing selection is often documented for inversion polymorphisms
448 (Wellenreuther & Bernatchez 2018). Inversions can extend the impact of barrier loci (reproduc-
449 tive isolation loci) to linked loci and promote additional barrier loci accumulating in inverted
450 regions with gene flow between populations which facilitates the efficient spread of adaptive
451 standing variance in inversions (Morjan & Rieseberg 2004). After the establishment of local
452 LD, gene drift can facilitate disruptive selection (Dieckmann & Doebeli 1999) which suggests
453 random process may play a role in speciation.

454
455 There are some studies revealing that a large scale of SNPs for some key adaptive phenotypes
456 or environmental axes could be explained by SNPs in chromosome inversions (Morales et al.
457 2019). Westram also found that 75% of non-neutral SNPs cluster together in three LGs (6, 14,
458 17) which suggests that large regions of low recombination are consistent with the presence of
459 chromosomal rearrangements.

460 4.5 Future direction

461 From a speciation perspective, the main goal is to infer the number, distribution and intensity
462 of gene flow barriers, and their impact on other genomic regions, then find genetic speciation is-
463 land which harbor loci underlying reproductive isolation. However, due to many other processes'
464 existence, it is challenging to detect this signal from genome scan. Chromosome inversions are
465 regarded as candidates for gene flow barriers. Then the research on ages and spatial distribu-
466 tions of inversions will be important to understand local adaptation and evolution history of
467 reproductive isolation in *Littorina saxatilis*. It is premature to use diversity and divergence data
468 of *Littorina saxatilis* to calculate exact inversion ages, but Morals (Morales et al. 2019) inferred
469 from Sweden snails data that the origins of inversions are earlier than postglacial colonization
470 of the Swedish coast.

471
472 After finding inversions regions under divergence selection, the next step is to identify why they
473 are so differentiated. These regions are likely to contain reproductive isolation loci which need
474 more exploration, such as finding break point of inversions, functional annotation of outlier loci,
475 selection scans, phenotype-habitat-genotype associations, experimental analyses with sequencing
476 (Koch et al. 2021). It is worth mentioning that size is the feature under strong divergence

477 selection and also could be the main reason for assortative mating between ecotypes (Smadja &
478 Butlin 2011), therefore size is likely to facilitate the formation of barriers between ecotypes. Then
479 finding the association between size and reproductive isolation which is deduced by Johannesson
480 (Johannesson 2016) is another promising research direction.

481 **5 Conclusion**

482 In general, I found two distinct groups in the Spanish system of *Littorina saxatilis* instead of a
483 cline displayed in the more recent system Sweden. Compared with the wave ecotype, hybrids
484 are much closer to the crab ecotype. Besides, the directional introgression from wave to crab
485 destroys high LD in the crab ecotype. Differentiation analysis hints at balancing selection on
486 this species and LD analysis hints at stronger selection and more polymorphic inversions of crab
487 ecotype. Also, I found some new chromosome inversions candidates for future verification. This
488 research tried to infer the process and mechanisms of *Littorina saxatilis* ecological divergence and
489 speciation after the phase of the Sweden system. However, to explore the underlying mechanisms
490 of parallel ecological divergence with gene flow and the evolution of reproductive isolation with
491 interbreeding, more efforts are still needed in future.

492 **Code availability**

493 All of the scripts, software parameter settings, and dependencies packages are included in the
494 following website: <https://github.com/rz520/LittorinaPipeline>

495 Acknowledgement

496 Here I would like to appreciate supervisor Matteo Fumagalli for his continuous help on software
497 usage, parameter setting, and results interpretation. Co-supervisor Francesca Raffini and Roger
498 K. Butlin gave me lots of guidance on *Littorina saxatilis* background knowledge and recent
499 research progress. All of my supervisors helped me to polish my thesis by giving constructive
500 suggestion. I am also grateful to other master students in Matteo's Lab and my friends for their
501 help, encouragement and company.

502 **References**

- 503 Beaumont, M. A. (2010), ‘Approximate bayesian computation in evolution and ecology’, *Annual
504 Review of Ecology, Evolution, and Systematics* **41**(1), 379–406.
505 **URL:** <https://doi.org/10.1146/annurev-ecolsys-102209-144621>
- 506 Butlin, R. K., Galindo, J. & Grahame, J. W. (2008), ‘Sympatric, parapatric or allopatric: the
507 most important way to classify speciation?’, *Philosophical Transactions of the Royal Society
508 B: Biological Sciences* **363**(1506), 2997–3007.
509 **URL:** <https://royalsocietypublishing.org/doi/abs/10.1098/rstb.2008.0076>
- 510 Cruickshank, T. E. & Hahn, M. W. (2014), ‘Reanalysis suggests that genomic islands of specia-
511 tion are due to reduced diversity, not reduced gene flow’, *Molecular Ecology* **23**(13), 3133–3157.
512 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.12796>
- 513 Dieckmann, U. & Doebeli, M. (1999), On the origin of species by sympatric speciation, Iiasa
514 interim report, IIASA, Laxenburg, Austria.
515 **URL:** <http://pure.iiasa.ac.at/id/eprint/5926/>
- 516 Faria, R., Chaube, P., Morales, H. E., Larsson, T., Lemmon, A. R., Lemmon, E. M., Rafajlović,
517 M., Panova, M., Ravinet, M., Johannesson, K., Westram, A. M. & Butlin, R. K. (2019),
518 ‘Multiple chromosomal rearrangements in a hybrid zone between littorina saxatilis ecotypes’,
519 *Molecular Ecology* **28**(6), 1375–1393.
520 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.14972>
- 521 Fitzpatrick, B. M., Fordyce, J. A. & Gavrilets, S. (2008), ‘What, if anything, is sympatric
522 speciation?’, *Journal of Evolutionary Biology* **21**(6), 1452–1459.
523 **URL:** <https://doi.org/10.1111/j.1420-9101.2008.01611.x>
- 524 Fox, E. A., Wright, A. E., Fumagalli, M. & Vieira, F. G. (2019), ‘ngsLD: evaluating linkage
525 disequilibrium using genotype likelihoods’, *Bioinformatics* **35**(19), 3855–3856.
526 **URL:** <https://doi.org/10.1093/bioinformatics/btz200>
- 527 Galindo, J., Martínez-Fernández, M., Rodríguez-Ramilo, S. T. & Rolán-Alvarez, E. (2013),
528 ‘The role of local ecology during hybridization at the initial stages of ecological speciation in
529 a marine snail’, *Journal of Evolutionary Biology* **26**(7), 1472–1487.
530 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/jeb.12152>
- 531 Galindo, J., Morán, P. & Rolán-Alvarez, E. (2009), ‘Comparing geographical genetic differentia-
532 tion between candidate and noncandidate loci for adaptation strengthens support for parallel
533 ecological divergence in the marine snail littorina saxatilis’, *Molecular Ecology* **18**(5), 919–930.
534 **URL:** <https://doi.org/10.1111/j.1365-294X.2008.04076.x>
- 535 Gilly, A., Southam, L., Suveges, D., Kuchenbaecker, K., Moore, R., Melloni, G. E. M., Hatziko-
536 toulas, K., Farmaki, A.-E., Ritchie, G., Schwartzentuber, J., Danecek, P., Kilian, B., Pollard,

- 537 M. O., Ge, X., Tsafantakis, E., Dedoussis, G. & Zeggini, E. (2018), 'Very low-depth whole-
538 genome sequencing in complex trait association studies', *Bioinformatics* **35**(15), 2555–2561.
539 **URL:** <https://doi.org/10.1093/bioinformatics/bty1032>
- 540 Grahame, J. W., Wilding, C. S. & Butlin, R. K. (2006), 'Adaptation to a steep environ-
541 mental gradient and an associated barrier to gene exchange in littorina saxatilis', *Evolution*
542 **60**(2), 268–278.
543 **URL:** <https://doi.org/10.1111/j.0014-3820.2006.tb01105.x>
- 544 Hollander, J., Lindegarth, M. & Johannesson, K. (2005), 'Local adaptation but not geographical
545 separation promotes assortative mating in a snail', *Animal Behaviour* **70**(5), 1209–1219.
546 **URL:** <https://www.sciencedirect.com/science/article/pii/S0003347205002733>
- 547 Hudson, R. R., Slatkin, M. & Maddison, W. P. (1992), 'Estimation of levels of gene flow from
548 dna sequence data.', *Genetics* **132**(2), 583–589.
549 **URL:** <https://www.genetics.org/content/132/2/583>
- 550 Johannesson, K. (2016), 'What can be learnt from a snail?', *Evolutionary Applications* **9**(1), 153–
551 165.
552 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/eva.12277>
- 553 Johannesson, K., Johannesson, B. & Rolán-Alvarez, E. (1993), 'Morphological differentiation
554 and genetic cohesiveness over a microenvironmental gradient in the marine snail littorina
555 saxatilis', *Evolution* **47**(6), 1770–1787.
556 **URL:** <https://doi.org/10.1111/j.1558-5646.1993.tb01268.x>
- 557 Johannesson, K., Panova, M., Kemppainen, P., André, C., Rolán-Alvarez, E. & Butlin, R. K.
558 (2010), 'Repeated evolution of reproductive isolation in a marine snail: unveiling mecha-
559 nisms of speciation', *Philosophical Transactions of the Royal Society B: Biological Sciences*
560 **365**(1547), 1735–1747.
561 **URL:** <https://royalsocietypublishing.org/doi/abs/10.1098/rstb.2009.0256>
- 562 Johannesson, K., Rolán-Alvarez, E. & Ekendahl, A. (1995), 'Incipient reproductive isolation
563 between two sympatric morphs of the intertidal snail littorina saxatilis', *Evolution* **49**(6), 1180–
564 1190.
565 **URL:** <https://doi.org/10.1111/j.1558-5646.1995.tb04445.x>
- 566 Kirkpatrick, M. & Ravigné, V. (2002), 'Speciation by natural and sexual selection: Models and
567 experiments', *The American Naturalist* **159**(S3), S22–S35.
568 **URL:** <http://www.jstor.org/stable/10.1086/338370>
- 569 Koch, A., Brierley, C. & Lewis, S. L. (2021), 'Effects of earth system feedbacks on the potential
570 mitigation of large-scale tropical forest restoration', *Biogeosciences* **18**(8), 2627–2647.
571 **URL:** <https://bg.copernicus.org/articles/18/2627/2021/>
- 572 Korneliussen, T., Moltke, I., Albrechtsen, A. & Nielsen, R. (2013), 'Calculation of tajima's d
573 and other neutrality test statistics from low depth next-generation sequencing data', *BMC*

- 574 *Bioinformatics* **14**(1), 289.
575 **URL:** <http://www.biomedcentral.com/1471-2105/14/289>
- 576 Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. (2014), ‘ANGSD: Analysis of next generation
577 sequencing data’, *BMC Bioinformatics* **15**(1), 356.
578 **URL:** <http://www.biomedcentral.com/1471-2105/15/356/abstract>
- 579 Meisner, J. & Albrechtsen, A. (2018), ‘Inferring population structure and admixture proportions
580 in low-depth ngs data’, *Genetics* **210**(2), 719–731.
581 **URL:** <https://www.genetics.org/content/210/2/719>
- 582 Morales, H. E., Faria, R., Johannesson, K., Larsson, T., Panova, M., Westram, A. M. & Butlin,
583 R. K. (2019), ‘Genomic architecture of parallel ecological divergence: Beyond a single envi-
584 ronmental contrast’, *Science Advances* **5**(12).
585 **URL:** <https://advances.sciencemag.org/content/5/12/eaav9963>
- 586 Morjan, C. L. & Rieseberg, L. H. (2004), ‘How species evolve collectively: implications of gene
587 flow and selection for the spread of advantageous alleles’, *Molecular Ecology* **13**(6), 1341–1356.
588 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1365-294X.2004.02164.x>
- 589 Panova, M., Aronsson, H., Cameron, R., Dahl, P., Godhe, A., Lind, U., Ortega-Martinez,
590 O., Pereyra, R., Tesson, S., Wrangé, A.-L., Blomberg, A. & Johannesson, K. (2016), *DNA
591 Extraction Protocols for Whole-Genome Sequencing in Marine Organisms*, Vol. 1452, pp. 13–
592 44.
593 **URL:** <https://pubmed.ncbi.nlm.nih.gov/27460368/>
- 594 Panova, M., Blakeslee, A. M. H., Miller, A. W., Mäkinen, T., Ruiz, G. M., Johannesson, K.
595 & André, C. (2011), ‘Glacial history of the north atlantic marine snail, littorina saxatilis,
596 inferred from distribution of mitochondrial dna lineages’, *PLOS ONE* **6**(3), 1–14.
597 **URL:** <https://doi.org/10.1371/journal.pone.0017511>
- 598 Panova, M., Hollander, J. & Johannesson, K. (2006), ‘Site-specific genetic divergence in parallel
599 hybrid zones suggests nonallopatric evolution of reproductive barriers’, *Molecular Ecology*
600 **15**(13), 4021–4031.
601 **URL:** <https://doi.org/10.1111/j.1365-294X.2006.03067.x>
- 602 Rafajlovic, M., Eriksson, A., Rimark, A., Hintz-Saltin, S., Charrier, G., Panova, M., Andre, C.,
603 Johannesson, K. & Mehlig, B. (2013), ‘The Effect of Multiple Paternity on Genetic Diversity
604 of Small Populations during and after Colonisation’, *PLOS ONE* **8**(10).
605 **URL:** <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0075587>
- 606 Ravinet, M., Faria, R., Butlin, R. K., Galindo, J., Bierne, N., Rafajlović, M., Noor, M. A. F.,
607 Mehlig, B. & Westram, A. M. (2017), ‘Interpreting the genomic landscape of speciation: a
608 road map for finding barriers to gene flow’, *Journal of Evolutionary Biology* **30**(8), 1450–1477.
609 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/jeb.13047>

- 610 Ravinet, M., Westram, A., Johannesson, K., Butlin, R., André, C. & Panova, M. (2016), 'Shared
611 and nonshared genomic divergence in parallel ecotypes of *Littorina saxatilis* at a local scale',
612 *Molecular Ecology* **25**(1), 287–305.
613 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.13332>
- 614 Reid, D. G. (1996), *Systematics and evolution of Littorina*, Ray Society.
615 **URL:** <https://www.nhbs.com/systematics-and-evolution-of-littorina-book>
- 616 Reynolds, J., Weir, B. S. & Cockerham, C. C. (1983), 'Estimation of the coancestry coefficient:
617 Basis for a short-term genetic distance', *Genetics* **105**(3), 767–779.
618 **URL:** <https://www.genetics.org/content/105/3/767>
- 619 Rolan-Alvarez, E., Johannesson, K. & Erlandsson, J. (1997), 'The maintenance of a cline in
620 the marine snail *Littorina saxatilis*: The role of home site advantage and hybrid fitness',
621 *EVOLUTION* **51**(6), 1838–1847.
622 **URL:** <https://doi.org/10.1111/j.1558-5646.1997.tb05107.x>
- 623 Schlüter, D. (2009), 'Evidence for ecological speciation and its alternative', *Science*
624 **323**(5915), 737–741.
625 **URL:** <https://science.sciencemag.org/content/323/5915/737>
- 626 Skotte, L., Korneliussen, T. S. & Albrechtsen, A. (2013), 'Estimating Individual Admixture
627 Proportions from Next Generation Sequencing Data', *Genetics* **195**(3), 693–702.
628 **URL:** <https://doi.org/10.1534/genetics.113.154138>
- 629 Smadja, C. M. & Butlin, R. K. (2011), 'A framework for comparing processes of speciation in
630 the presence of gene flow', *MOLECULAR ECOLOGY* **20**(24), 5123–5140.
631 **URL:** <https://doi.org/10.1111/j.1365-294X.2011.05350.x>
- 632 Therkildsen, N. O., Wilder, A. P., Conover, D. O., Munch, S. B., Baumann, H. & Palumbi,
633 S. R. (2019), 'Contrasting genomic shifts underlie parallel phenotypic evolution in response
634 to fishing', *Science* **365**(6452), 487–490.
635 **URL:** <https://science.sciencemag.org/content/365/6452/487>
- 636 Wellenreuther, M. & Bernatchez, L. (2018), 'Eco-evolutionary genomics of chromosomal inver-
637 sions', *Trends in Ecology & Evolution* **33**(6), 427–440.
638 **URL:** <https://www.sciencedirect.com/science/article/pii/S0169534718300788>
- 639 Westram, A. M., Galindo, J., Alm Rosenblad, M., Grahame, J. W., Panova, M. & Butlin,
640 R. K. (2014), 'Do the same genes underlie parallel phenotypic divergence in different *Littorina
641 saxatilis* populations?', *Molecular Ecology* **23**(18), 4603–4616.
642 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.12883>
- 643 Westram, A. M., Rafajlović, M., Chaube, P., Faria, R., Larsson, T., Panova, M., Ravinet, M.,
644 Blomberg, A., Mehlig, B., Johannesson, K. & Butlin, R. (2018), 'Clines on the seashore: The
645 genomic architecture underlying rapid divergence in the face of gene flow', *Evolution Letters*

- 646 **2**(4), 297–309.
- 647 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1002/evl3.74>
- 648 Wilder, A. P., Palumbi, S. R., Conover, D. O. & Therkildsen, N. O. (2020), ‘Footprints of local
649 adaptation span hundreds of linked genes in the atlantic silverside genome’, *Evolution Letters*
650 **4**(5), 430–443.
- 651 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1002/evl3.189>
- 652 Yeaman, S. (2013), ‘Genomic rearrangements and the evolution of clusters of locally adaptive
653 loci’, *Proceedings of the National Academy of Sciences* **110**(19), E1743–E1751.
- 654 **URL:** <https://www.pnas.org/content/110/19/E1743>
- 655 Zhang, C., Dong, S.-S., Xu, J.-Y., He, W.-M. & Yang, T.-L. (2018), ‘PopLDdecay: a fast
656 and effective tool for linkage disequilibrium decay analysis based on variant call format files’,
657 *Bioinformatics* **35**(10), 1786–1788.
- 658 **URL:** <https://doi.org/10.1093/bioinformatics/bty875>

659 **Supplementary information**

660 **A Data collection and processing**

661 **A.1 Sampling**

662 *Littorina saxatilis* snails were sampled in March 2017 on the west coast of Galicia in Spain:
663 Centinela (hereafter ER_EA1; N 42° 4' 38.06", W 8° 53' 47.47"). The sample consisted of
664 approximately 600 snails from two 5 m wide bands and separated by 0-5 m, stretching from the
665 upper limit of the species distribution in the splash zone to its lower limit close to low water of
666 spring tides. Sampling was denser in the lower part of the shore where hybridization between
667 2 ecotypes described before was expected (Galindo et al. 2013). Snails were sampled by hand,
668 haphazardly, as far as possible without reference to phenotype but aiming to avoid juveniles.
669 All of snails were stored in individual tube, moistened with seawater and kept at 4 °C until
670 phenotyping was complete, then the head and foot of each snail was preserved in 100% ethanol.

671 **A.2 DNA sequencing and processing**

672 Low coverage whole genome sequencing (lcWGS) was conducted for 73 snails from ER_EA1,
673 chosen to cover the full sampling range. DNA was extracted using the modified CTAB protocol
674 described by Panova (Panova et al. 2016). Library preparation (in-house high-throughput gDNA
675 library prep) and sequencing (Illumina HiSeq4000, 150bp PE) were carried out by The Oxford
676 Genomics Centre with target coverage 3x based on the estimated genome size of 1.35Gb. Samples
677 were sequenced in eight lanes, for a total of 16 lanes for each individual. Raw reads were trimmed
678 with Trimmomatic v.0.38 (<https://github.com/usadellab/Trimmomatic>) to remove the Illumina
679 adapters, mapped to the *L. saxatilis* draft reference genome (Swedend, crab; 1.35Gb) (Westram
680 et al. 2018) using bwa mem v.0.7.17 (<https://github.com/lh3/bwa>) and default settings, and fil-
681 tered by base and map quality 20 with Samtools v.1.7 (<https://github.com/samtools/samtools>).
682 PCR duplicates and PE overlap were removed with Picard v.2.0.1 ([https://github.com/broadin-
683 titute/picard](https://github.com/broadinstitute/picard)) and bamUtil v.1.0.15 (<http://genome.sph.umich.edu/wiki/BamUtil>), respectively.
684 Bam files belonging to the same individual were sorted and merged with Samtools.

685 **B Table and plots**

686 **B.1 Mean summary statistics**

Linkage group (LG)	LD cluster	Cluster size (cM)	Start (cM)	End (cM)	Mean FST	Mean PBS	Mean Pi	Mean LD (r2)	candidate
LG1	LGC1.1	2.1	0	2.1	0.0019398	0.012677	0.00648093	0.1939411	1
LG1	LGC1.2	5.42	75.53	80.95	0.00219029	0.02014028	0.00536197	0.2134792	1
LG1	LG1	80.95	0	80.95	0.00212763	0.01073747	0.00594011	0.1704242	0
LG2	LGC2.1	13.87	0.34	14.21	0.00614239	0.01879701	0.00576577	0.1999839	1
LG2	LG2	88.76	0	88.76	0.00274161	0.01610663	0.00592564	0.1646505	0
LG3	LG3	68.02	0	68.02	0.00288605	0.01137463	0.00603668	0.1761798	0
LG4	LGC4.1	0.48	1.03	1.51	0.00463436	0.03770791	0.00587387	0.1847306	1
LG4	LG4	56.52	0	56.52	0.00243	0.0140493	0.00587649	0.1711099	0
LG5	LG5	54.07	0	54.07	0.00280035	0.01904912	0.00591619	0.1699921	0
LG6	LGC6.1	29.3	0	29.3	0.00297457	0.03524753	0.00567617	0.1820881	1
LG6	LGC6.2	20.57	8.73	29.3	0.00269045	0.02541398	0.00579544	0.1825618	1
LG6	LG6	60.25	0	60.25	0.00227857	0.01616375	0.00634677	0.1691389	0
LG7	LGC7.1	1.73	36.01	37.74	0.00275902	0.01656839	0.00683345	0.2117835	1
LG7	LGC7.2	9.29	42.08	51.37	0.00364395	0.02627485	0.00641327	0.1799024	1
LG7	LG7	51.37	0	51.37	0.00235047	0.0124258	0.00629948	0.1714745	0
LG8	LG8	54.38	0	54.38	0.0025284	0.01993345	0.00671042	0.186891	0
LG9	LGC9.1	23.18	18.64	41.82	0.00275828	0.02324772	0.00566658	0.1917811	1
LG9	LG9	56.67	0	56.67	0.00256584	0.02090745	0.00763238	0.1833917	0
LG10	LGC10.1	2.54	0.58	3.12	0.00251415	0.01067118	0.00638226	0.1959969	1
LG10	LG10	45.53	0	45.53	0.00280773	0.02139548	0.00636286	0.188335	0
LG11	LGC11.1	0.59	52.32	52.91	0.00259493	0.0145582	0.00686731	0.2100456	1
LG11	LG11	58.39	0	58.39	0.0024994	0.02136298	0.00675471	0.1950033	0
LG12	LGC12.1	26.31	3.32	29.63	0.00155347	0.01451515	0.00667943	0.1775293	1
LG12	LGC12.2	11.52	48.71	60.24	0.00299152	0.01851379	0.00654778	0.1870268	1
LG12	LG12	60.24	0	60.24	0.00314249	0.02937913	0.00609529	0.1726986	0
LG13	LG13	51.3	0	51.3	0.0018803	0.01141404	0.00618406	0.1804347	0
LG14	LGC14.1	11.32	0.39	11.71	0.00325176	0.00927283	0.0060404	0.1966089	1
LG14	LGC14.2	2.9	8.81	11.71	0.00286758	0.00959897	0.00606569	0.1960541	1
LG14	LGC14.3	23.23	11.71	34.94	0.0034732	0.01206736	0.00528926	0.1782203	1
LG14	LG14	48.71	0	48.71	0.00332075	0.00996925	0.00583716	0.1881616	0
LG15	LG15	56.49	0	56.49	0.00201495	0.00941298	0.00633721	0.1755119	0
LG16	LG16	57.44	0	57.44	0.00212438	0.01346309	0.00695419	0.1911745	0
LG17	LGC17.1	15.33	46.99	62.32	0.00561686	0.0379331	0.00655994	0.1748176	1
LG17	LG17	62.32	0	62.32	0.00219622	0.01432516	0.00598865	0.1672243	0

Table 1: Mean summary statistics in inversion candidates and whole chromosome

687 B.2 PBS

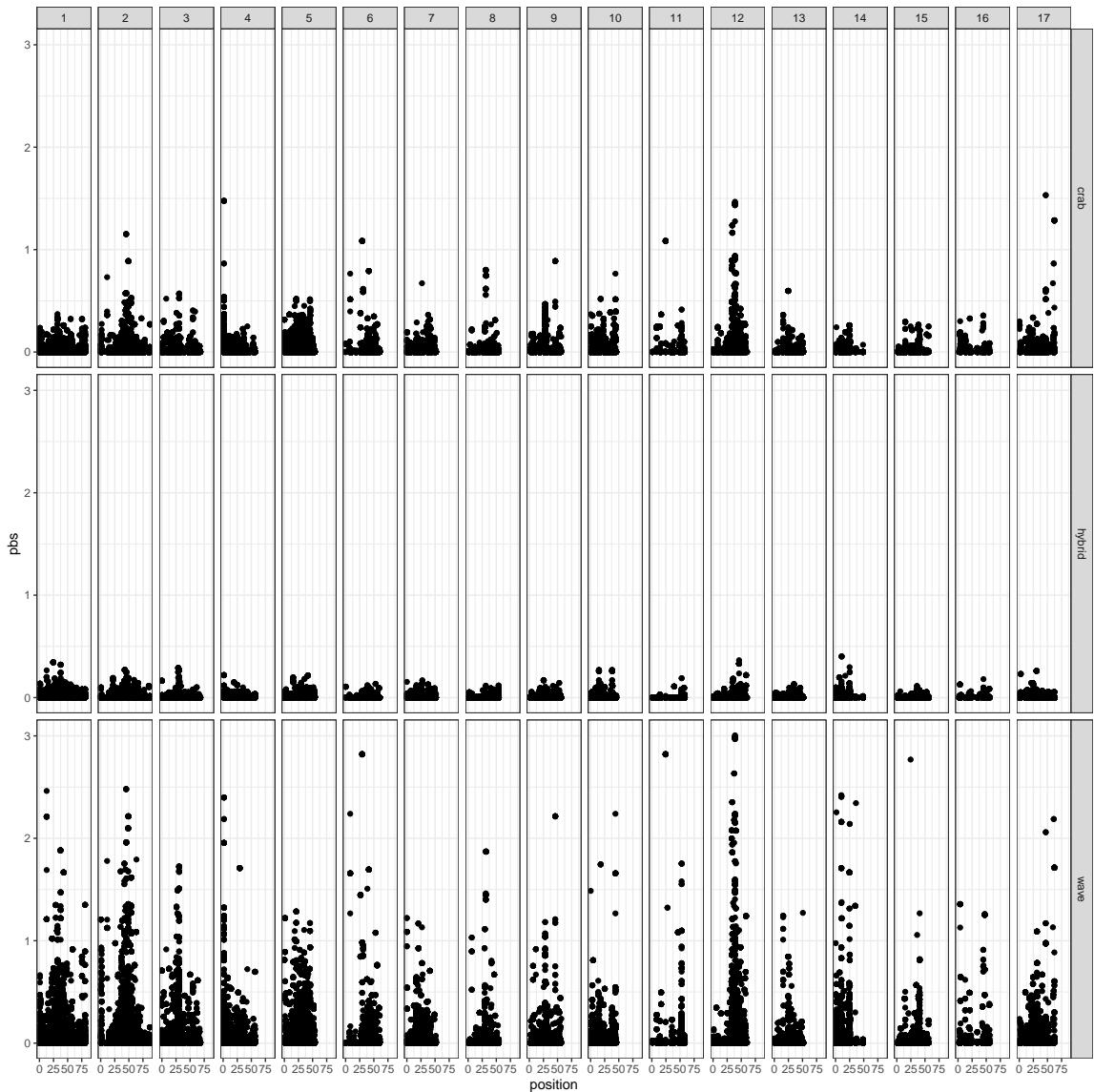


Figure 1: PBS value scatter plot of every chromosome and population in windows

688 B.3 Pi

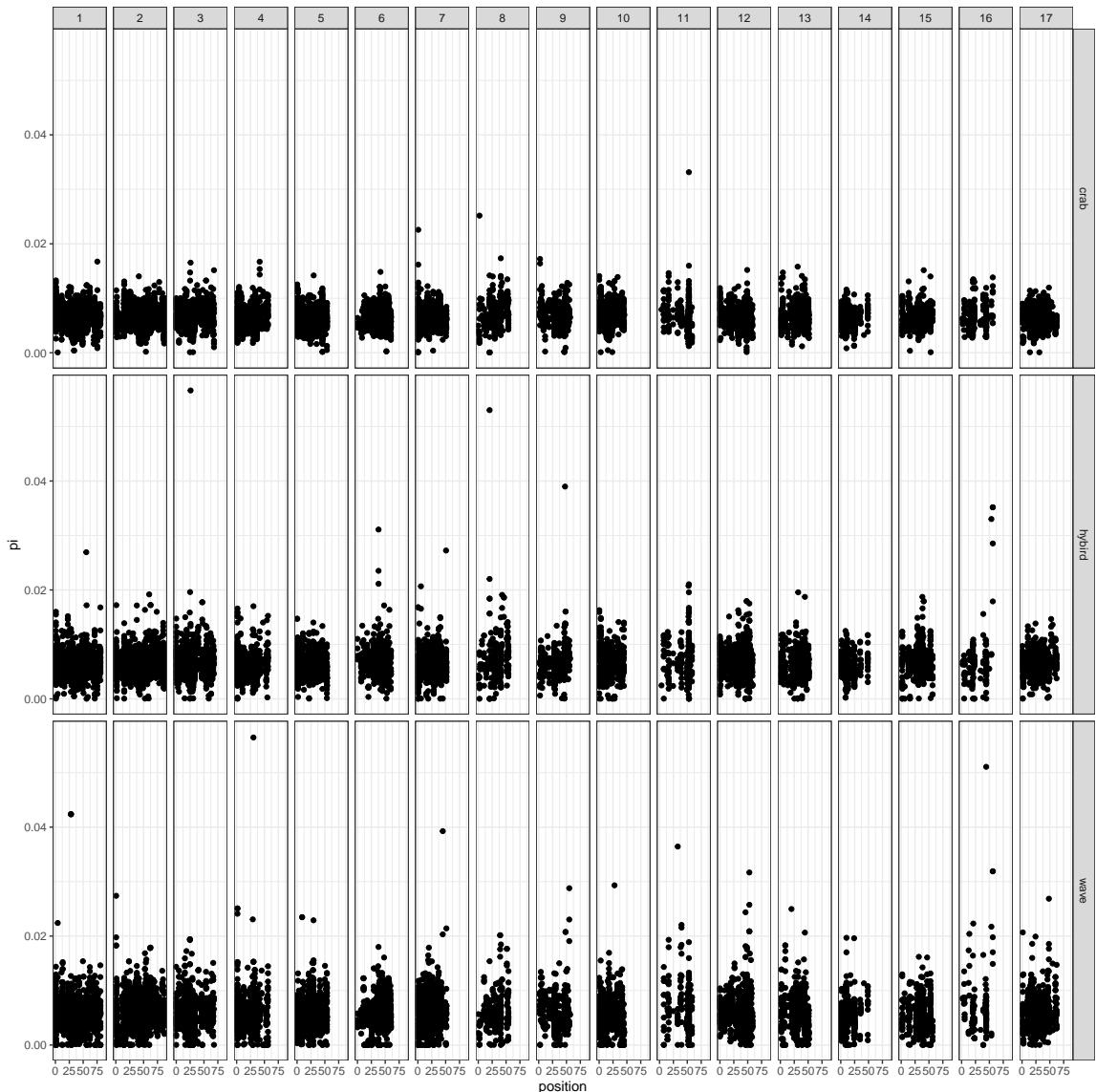


Figure 2: Pi value scatter plot of every chromosome and population in windows

689 **B.4 LD decay**

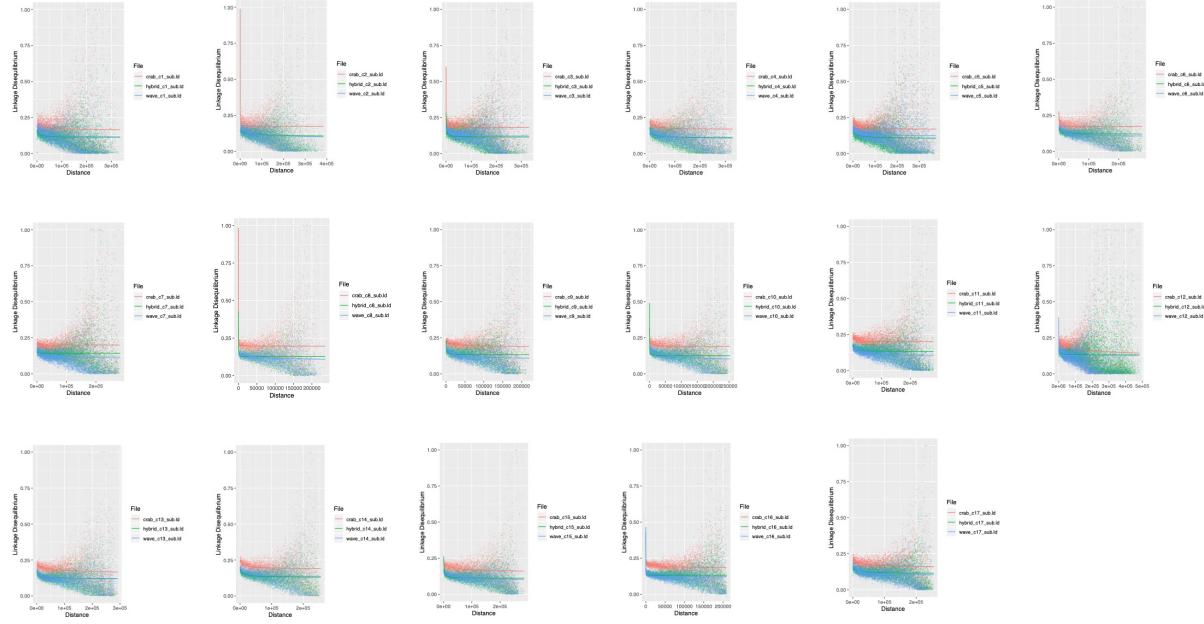


Figure 3: LD decay analysis of every chromosome

690 **B.5 LD heat map**

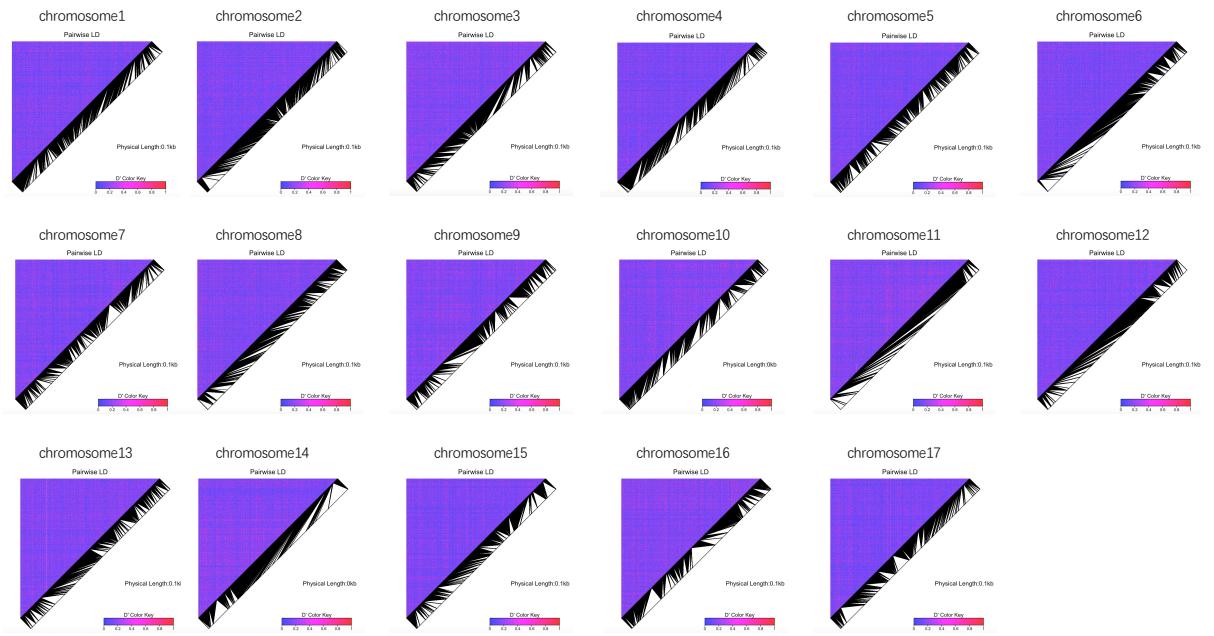


Figure 4: LD heat map of every chromosome